

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054

---

## Text-to-Image: The Future Generation

---

### Abstract

This study conducts qualitative and quantitative comparisons among four state-of-the-art text-to-image generative models: Stable diffusion, DALLE, CLIP, and Kandinsky. The focus is on evaluating image quality, diversity, and alignment with textual descriptions, highlighting the superiority of Controllable Generative AI. Additionally, to tackle the challenge of aligning creativity with functionality, we propose a framework for Controllable Generative AI in Biomimetic Design, integrating constraints and control over the creative process to generate innovative and feasible designs. The code for our implementation is accessible via the link: <https://github.com/JohnsonForSure/ML-in-Bioinfo-and-Healthcare>.

### 1. Introduction

Our perception of the world encompasses multiple sensory modalities, including sight, sound, touch, smell, and taste (Baltrušaitis et al., 2018). Multimodal machine learning endeavors to develop models capable of interpreting and associating information from these diverse modalities (Yadav & Vishwakarma, 2020). Among the various multimodal tasks, text-to-image (T2I) synthesis stands out, involving the generation of images corresponding to given textual descriptions. This task holds immense potential across various domains such as image editing, video games, and computer-aided design, necessitating a deep comprehension of the entities being created (Ramesh et al., 2022; Liu et al., 2021).

Recent advancements in Generative Adversarial Networks (GANs), autoregressive models, and diffusion models have propelled the frontier of T2I synthesis (Zameshina et al., 2023). GANs engage in a competitive framework with a generator and discriminator network, striving to produce realistic images indistinguishable from genuine data (Van Looveren et al., 2021). Autoregressive models, conversely, generate data iteratively, conditioning each step on preceding information (Du et al., 2022). These models, based on transformer architectures, utilize self-attention mechanisms to capture dependencies and relationships within input sequences (Esser et al., 2021). On the other hand, diffusion models represent a paradigm shift in machine learning, in-

troducing noise to data and progressively refining it through denoising iterations (Meng et al., 2021).

To evaluate the efficacy of these generative AI models, we conduct qualitative and quantitative comparisons across four state-of-the-art T2I generation paradigms: Stable diffusion (Rombach et al., 2022), DALLE (Ramesh et al., 2021), CLIP (Ramesh et al., 2022), and Kandinsky generative AI model (Razzhigaev et al., 2023). Our assessment examines the quality, diversity, and alignment of generated images with given textual descriptions across various levels, with a focus on highlighting the superior performance of Generative AI in text-to-image (T2I) synthesis. While these generative AI models showcase remarkable creativity, their unrestricted nature may not consistently align with practical applications, particularly in biomimetic design.

Biomimicry involves replicating natural structures and processes to address complex challenges (Vincent et al., 2006). Meanwhile, the design space must be bounded by factors such as living structures, movement mechanisms, and scientific principles. To address this challenge, we propose a framework for Controllable Generative AI in Biomimetic Design.

Our proposed framework utilizes Large Language Models (LLMs) and morphology-constrained text-to-image (T2I) transformation. It begins with the generation of a table of entities associated with specific adjectives describing the target's nature. These entities are then projected into a unified vector space, with outliers removed based on similarity. Subsequently, adjectives are collected, ranked, and filtered to eliminate conflicting attributes. An LLM selects the best-aligned entity, and its morphological features guide the T2I transformation. This approach ensures that the generated designs align with functional requirements and real-world constraints, making them suitable for practical applications in biomimicry and related fields. By bridging the gap between AI's boundless creativity and the necessity for implementable solutions, our framework, integrating ControlNet (Zhang et al., 2023) and LLMs, holds the potential to significantly advance the realm of biomimetic design, ushering in new possibilities for innovative and sustainable solutions inspired by nature.

## 055 2. Literature Review

### 056 2.1. Stable Diffusion model ([Rombach et al., 2022](#))

057  
058 The Diffusion model differs from the other generation methods in that it needs to be run in multiple steps in terms  
059 of timing and the hidden variable  $z$  has the same dimension  
060 as the original size of  $x$ . For Stable diffusion model.  
061 The techniques involved in Stable Diffusion are Diffusion  
062 Model (DDPM), Attention, Autoencoder. It based on Latent  
063 Diffusion Models (LDMs)

064  
065 Compared to the original diffusion, the advantage is that the  
066 encoder is used to map the image from the original pixel  
067 space to the latent space, which reduces the amount of com-  
068 putation. Since diffusion model favors semantic compres-  
069 sion (which is easy to be perceived by the human eye) and  
070 GAN/AutoEncoder favors perceptual compression (which  
071 means that it is difficult for the human eye to perceive the  
072 high-frequency information), the image can be similarly  
073 low-pass filtered with the encoder first. Encoder can ef-  
074 fectively reduce the image edge length to 1/4 to 1/16 of  
075 the original, and even improve the quality of the generated  
076 samples.

### 077 2.2. Autoregressive model: CLIP and DALLE (both 078 from OpenAI)

#### 079 2.2.1. CLIP: CONTRASTIVE LANGUAGE-IMAGE 080 PRE-TRAINING ([RADFORD ET AL., 2021](#); 081 [RAMESH ET AL., 2022](#))

082 The approach of CLIP centers on leveraging a large corpus  
083 of image-text pairs to enforce alignment within a shared  
084 feature space. In this framework, a 12-layer transformer  
085 is employed for text encoding, while image encoding can  
086 utilize either ResNet or Vision Transformer (ViT), with ViT  
087 generally outperforming ResNet in evaluations.

088 The process is simplified to include linear projection of en-  
089 coder representations to a multi-modal embedding space,  
090 omitting non-linear projection. Symmetric cross-entropy  
091 loss is employed alongside the ADAM optimizer, with train-  
092 ing initiated from scratch. Notably, the task is defined as  
093 text-to-image (T2I) generation rather than exact word-to-  
094 image mapping in a single step, as the latter can be challeng-  
095 ing to train. Each batch consists of  $N$  texts and  $N$  images,  
096 with the model generating  $N^2$  scores.

097 Regarding data, CLIP relies on a dataset comprising approx-  
098 imately 400 million Text-Image pairs known as WebImage-  
099 Text, akin in scale to the WebText dataset used in GPT-2.  
100 Data augmentation is limited to random square crops from  
101 resized images.

102 Despite its strengths, CLIP exhibits limitations in special-  
103 ized domains such as EuroSAT/MNIST, and abstract tasks

104 like predicting the number or positional distance of items  
105 within an image, where zero-shot performance remains con-  
106 strained. Overcoming these limitations may necessitate en-  
107 hanced prompt engineering, sometimes involving iterative  
108 trial-and-error approaches.

#### 109 2.2.2. DALLE: LANGUAGE-IMAGE PRE-TRAINING 110 ([RAMESH ET AL., 2021](#))

111 In the DALL-E model, images undergo encoding via the  
112 VQGAN encoder, translating them into a sequence of tokens.  
113 Simultaneously, language descriptions are encoded using  
114 the BART encoder. These encoded descriptions and images  
115 are then fed into the BART decoder, an autoregressive model  
116 that predicts the next image. The loss function employed  
117 is softmax cross-entropy, computed between the model's  
118 prediction logits and the actual image encodings derived  
119 from the VQGAN.

120 During the inference stage, the caption is encoded by the  
121 BART encoder, and the beginning-of-sequence token is in-  
122 put to the BART decoder. Image tokens are sequentially  
123 sampled based on the decoder's predicted distribution over  
124 the next token, and sequences of image tokens are subse-  
125 quently decoded by the VQGAN decoder. CLIP is then  
126 utilized to select the best-generated images from the output..

#### 127 2.3. Kandinsky T2I Synthesis with Image Prior and 128 Latent Diffusion ([Razzhigaev et al., 2023](#))

129 The Kandinsky model represents an enhanced iteration in-  
130 corporating latent diffusion and image prior techniques. Nota-  
131 bly, it integrates CLIP-image embeddings instead of stan-  
132 dalone text encoders, improving image quality. Additionally,  
133 it adopts an image-prior approach, leveraging diffusion and  
134 linear mappings between CLIP's text and image embedding  
135 spaces while retaining additional conditioning with XLMR  
136 text embeddings. This dual-encoder setup encompasses  
137 CLIP-text with image prior mapping and XLMR, which  
138 remain frozen during the training phase.

139 A pivotal consideration driving the model's design was the  
140 efficiency observed in training latent diffusion models, as  
141 highlighted in Rombach et al. (2022). The model operates  
142 through three stages: text encoding, embedding mapping  
143 (image prior), and latent diffusion. A transformer-encoder  
144 model is employed at the embedding mapping stage, also  
145 referred to as the image prior. This component is trained  
146 from scratch using a diffusion process on text and image  
147 embeddings provided by the CLIP-ViT-L14 model. Notably,  
148 the model incorporates element-wise normalization of visual  
149 embeddings based on full-dataset statistics to expedite the  
150 convergence of the diffusion process. In the inference stage,  
151 inverse normalization is utilized to revert to the original  
152 CLIP-image embedding space.

### 110 3. Problem Statement

111 Text-to-image (T2I) synthesis, which generates visually  
112 realistic images from textual descriptions, stands at the  
113 crossroads of deep learning and computer vision, promising  
114 transformative possibilities. Recent strides in generative  
115 models like Generative Adversarial Networks (GANs),  
116 autoregressive models, and diffusion models have signifi-  
117 cantly expanded the horizons of this domain. Nevertheless,  
118 amidst these advancements, persistent challenges impede  
119 the seamless creation of high-quality, diverse, and contextu-  
120 ally aligned images from textual prompts.  
121

122 A fundamental hurdle lies in the diversity and fidelity of  
123 generative models themselves. GANs, renowned for their  
124 capacity to produce lifelike images, grapple with issues like  
125 mode collapse, training instability, and hyperparameter sen-  
126 sitivity (Sauer & Geiger, 2021; Li et al., 2023). Achieving  
127 the delicate equilibrium between generator and discrimi-  
128 nator networks demands meticulous parameter tuning to  
129 avert subpar outcomes. Furthermore, the adversarial dy-  
130 namics of GANs often hinder precise control over gener-  
131 ated imagery and the faithful portrayal of textual prompts  
132 (Kang et al., 2023). Autoregressive models, which construct  
133 data iteratively based on preceding information, struggle  
134 to ensure coherence and capture intricate textual nuances  
135 effectively (Esser et al., 2021). Their sequential nature can  
136 constrain their ability to grasp long-range dependencies  
137 and craft highly detailed visuals (Li et al., 2021). More-  
138 over, as input sequences lengthen, the computational de-  
139 mands of autoregressive models escalate, posing scalability  
140 challenges for generating high-resolution images (Du et al.,  
141 2022). While diffusion models present a promising avenue,  
142 they introduce their own complexities. The iterative process  
143 of adding and reversing noise demands significant compu-  
144 tational resources and meticulous optimization to maintain  
145 image quality throughout denoising iterations (Meng et al.,  
146 2021).

147 Beyond the idiosyncrasies of individual generative models,  
148 a broader challenge looms regarding these techniques' practical  
149 deployment and adaptability in real-world contexts.  
150 Integrating T2I synthesis into arenas such as image editing,  
151 gaming, and computer-aided design mandates tackling  
152 issues of scalability, interpretability, and user-centric cus-  
153 tomization (Wang et al., 2021). Ensuring generated images  
154 resonate with user preferences, exhibit consistency across  
155 diverse textual inputs, and seamlessly integrate into existing  
156 workflows poses an ongoing dilemma (Mayerson, 2023).

157 Addressing these challenges and propelling the T2I synthe-  
158 sis field forward necessitates comprehensive inquiries into  
159 their implications on generative model efficacy. Regrettably,  
160 comparative studies evaluating model performance across  
161 various challenges remain scarce in the current literature  
162 (Zhou & Shimada, 2023). While modern AI generators  
163

dazzle with their creativity, their unbounded imaginings often diverge from practical feasibility. In contexts like biomimicry or practical applications such as housing construction or biotechnology, unrestrained creativity risks impracticality. Constraints like structural integrity, biomechanics, airflow, and scientific principles demand a harmonious fusion of creativity with functionality. Effective creativity often hinges on seamlessly blending reality with imagination.

Thus, we present a framework to bridge the chasm between AI-generated creations and real-world applicability, particularly in biomimicry. By facilitating a synergy between AI generation and practical constraints, we aspire to unlock novel applications and visionary possibilities.

### 4. Methodology

#### 4.1. Generation Across Models

In our experiment, we deploy four models for text-to-image (T2I) tasks: the Stable diffusion, DALLE<sub>mini</sub>, CLIP, and Kandinsky models. Our approach entails generating images corresponding to a series of input sentences, each crafted to present a different level of complexity. These sentences range from simple depictions of everyday scenes to more intricate and abstract descriptions. Starting with straightforward imagery, the sentences gradually introduce nuances that challenge the model's ability to interpret and visualize the intended scenes. From descriptions of nature to celestial phenomena, the complexity of the vocabulary and concepts increases, pushing the boundaries of the model's understanding and creativity. Additionally, some sentences incorporate poetic language and unconventional scenarios, adding layers of complexity that require imaginative interpretation. Through this diverse set of input sentences, the method aims to explore the model's capacity to generate visually compelling images across a spectrum of linguistic complexity and conceptual abstraction. Furthermore, to extend the capabilities of controllable generative models, we deliberately select two sentences and compare the ControlNet's output with those from the four synthesis models, providing a comprehensive evaluation of performance with controllable constraints.

#### 4.2. Metric and Configurations

Evaluation of T2I model can be different. While, most of the machine evaluation process should be included in the training and testing process of the model which requires great load of memory and strong GPU. Here we only perform Human evaluation and IS Analysis to scrutinize the effectiveness of the models.

165 4.2.1. HUMAN EVALUATION

166 Human Evaluation focuses on confirming the credibility and  
 167 alignment of automatic evaluations through user analysis.  
 168 In our experiment, 8 students were recruited to conduct  
 169 the evaluation. Each Participant is asked to rate generated  
 170 images based on plausibility, including factors such as ob-  
 171 ject accuracy, counting, positional alignment, or image-text  
 172 alignment, as well as naturalness, assessing whether the  
 173 image appears realistic. Ratings are provided on a 5-Point  
 174 Likert scale, with 5 representing the best and 1 representing  
 175 the worst. Additionally, human evaluation is emphasized for  
 176 assessing challenging tasks and rare object combinations, as  
 177 it ensures accurate assessments and helps avoid bias related  
 178 to race or gender.

180 4.2.2. INCEPTION SCORE

182 Inception Score(IS) is a frequently used metric for analyzing  
 183 the performance of a generative model. A common form of  
 184 IS can be expressed as follows:

$$186 IS(G, z) = \exp(\mathbb{E}_{x \sim P(z)} D_{KL}(p(y|x) || p(y))) \quad (1)$$

187 Where  $D_{KL}()$  denotes the KL divergence distance between  
 188 two distributions, and  $p(y|x)$  and  $p(y)$  denotes the con-  
 189 ditional distribution and the marginal distribution.  $y$  is the  
 190 label prediction distribution given some classification model  
 191 (usually a deep neural network called the Inception V3 is  
 192 used), and  $x \sim P_n$  are images generated from a given gen-  
 193 erative model  $G$  sampled from the set of possible outcome  
 194 distribution  $P(z)$  given some text input vector  $z$ . In reality  
 195 we don't have the true distribution  $P(z)$  nor the conditional  
 196 and marginal distribution on  $y$ , hence we can only use a  
 197 emperical version of it. Let  $N$  be the total number of sam-  
 198 ples generated from the model  $G$  given a input vector  $z$ ,  
 199 and the set  $X_N(z) = \{x_1(z), \dots, x_N(z)\}$  to be the set of  
 200 generated images. Then we approximate  $p(y)$  with:  
 201

$$203 \hat{p}(y, z) = \frac{1}{N} \sum_{i=1}^N p_{classifier}(y|x_i(z)) \quad (2)$$

205 Here  $p_{classifier}(y|x_i(z))$  is the distribution output of some  
 206 chosen classifier given an input image  $x_i(z)$ . Then the  
 207 emperical version of the IS score can be written as:  
 208

$$210 \hat{IS}(G, z) = \exp\left(\frac{1}{N} \sum D_{KL}(p_{classifier}(y|x_i(z)) || \hat{p}(y, z))\right) \quad (3)$$

213 In our experiment, the classifier that we use is *placeholder*,  
 214 and we choose  $N = \text{placeholder}$  for all models.

216 4.3. Controllable Generative AI in Biomimetic Design

218 While modern AI generators can inspire with their creativity,  
 219 their limitless imaginings do not easily translate into reality.

Were such unconstrained creations employed in biomimicry or other practical applications like housing construction or biotechnology, they would likely prove infeasible. Factors such as living structures, movement mechanisms, airflow, geometry, and science constrain expression - creative flourishes must align with the functional. And the most impactful creativity is often the one integrating reality with imagination. Hence, Controlable generative AI is a promising advance, and by incorporating greater constraints and control over the creative process, it will open up even more opportunities in AI-powered fields in the future. On this point, here we provide a demo of how the Biomimetics process can benefit from the Large Language model (LLM) and the Controllable generative AI.

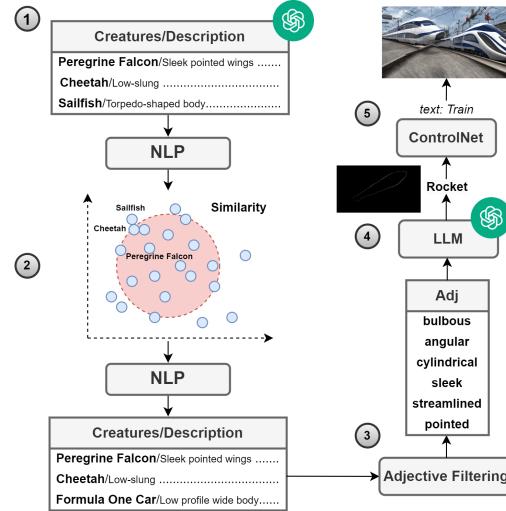


Figure 1. Framework for Integrating Controllable Generative AI into Biomimetic Design Processes.

Our framework comprises five steps, as shown in Figure 1:

- **Step1. Entity Generation and Selection**

The initial step in our approach involves generating a table of entities related to a specific adjective that describes the nature of the target. Two methods can be employed for this purpose. The first involves scraping existing data from online encyclopedias, such as Wikidata, or acquiring datasets from platforms like Kaggle or other databases. Alternatively, large language models (LLMs) can be directly accessed to assist in this process. LLMs, such as ChatGPT, which have made recent breakthroughs in Artificial Intelligence accessible to everyone, are powerful machine-learning models capable of understanding and generating natural language. These models are trained on extensive datasets of text and code, enabling them to learn patterns and relationships within the language. Utilizing interactive platforms, LLMs can generate a table of entities

and their descriptions that align with the specified requirements. It is important to note that the descriptions can encompass various aspects of the entities' nature, with morphology being the most straightforward to visualize.

#### • Step2. Entity Projection and Outlier Removal

A further processing step is crucial to focusing on the primary common features. This involves projecting all entities into a vector space based on their descriptions and performing dimension reduction to align them in a unified space. Subsequently, outliers exhibiting large discrepancies are identified and removed based on their similarity using Natural Language Processing (NLP) methods. In our implementation, we utilize the open-source Spacy library for our NLP tasks.

#### • Step3. Adjective Collection and Filtering

Following the morphological adjective similarity comparison, a bag of entities is obtained. All adjectives associated with these entities are collected and ranked based on their frequency of appearance. An additional step involves filtering out conflicting adjectives, such as "small" and "big" or "tall" and "short." If conflicting adjectives are present, both should be removed as they do not represent the target's primary influential nature.

#### • Step4. Entity Selection using LLMs

Utilizing the refined list of adjectives, an LLM is employed to select an entity that best aligns with all the adjectives. This process aims to identify an existing entity in the natural world that closely matches the desired characteristics. An online figure that best represents the target is then selected.

#### • Step5. Morphology-Constrained T2I Generation

Image analysis techniques are applied to the selected figure to extract its features. These features are then provided to the ControlNet as constraints to guide the T2I transformation process. By incorporating the constraints derived from the selected entity's morphological features, the generated images are expected to exhibit a fusion of the target image characteristics and the constraining features.

## 5. Experimental Results

In this section, we compare the generative results of four different methods: Stable Diffusion, DALL-E, CLIP, and Kandinsky. We will present some generated results of each method based on given text sentences in Section 5.1. Following that, we conduct qualitative and quantitative comparisons across these models in Section 5.2. Subsequently, we expand the constraints with ControlNet in Section 5.3. The

implementation details of the framework for Biomimetic design process integrating controllable generative AI and LLMs are provided in Appendix A.

### 5.1. Results Across Different Models

**Input Sentences:** The sentences were designed with varying levels of difficulty:

- *A red apple sits on a green leaf.* This sentence uses simple vocabulary, which is easy to understand because it depicts scenes from daily life.
- *An old oak tree stands tall in the forest, its branches reaching toward the sky.* Things get a little more complicated when you add descriptions that include adjectives like "old" and "tall." However, the concept of trees in the forest is still easy to imagine.
- *The fantastic glow of celestial bodies illuminates the canvas of the night, creating a mesmerizing tapestry of cosmic beauty.* This term uses more sophisticated vocabulary such as "celestial," "mesmerizing," and "tapestry." Thus, the resulting images also become more abstract and complex.
- *A firefly snakes its way through the velvety darkness of the forest.* More descriptive vocabulary, such as "snakes its way" and "velvety darkness." The model needs to generate a more poetic scene, rooted in nature.
- *The cat rides the horse.* Although simple vocabulary is used here, the concept of a cat riding a horse is so bizarre that it is difficult to deduce from modern image databases, ie., absurdity increases complexity.

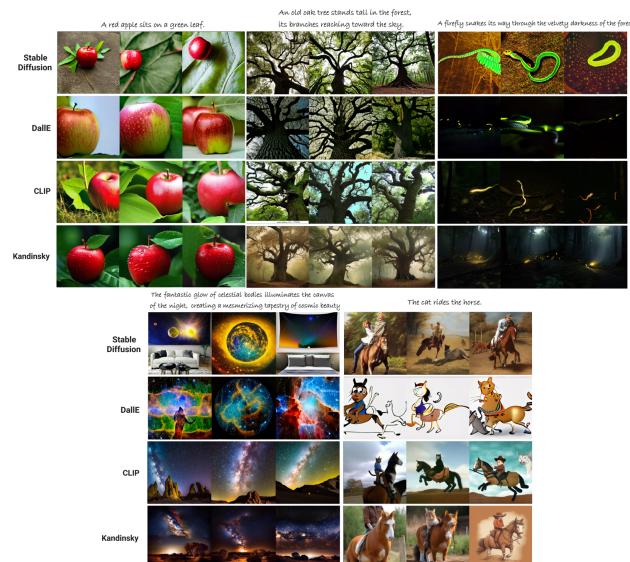
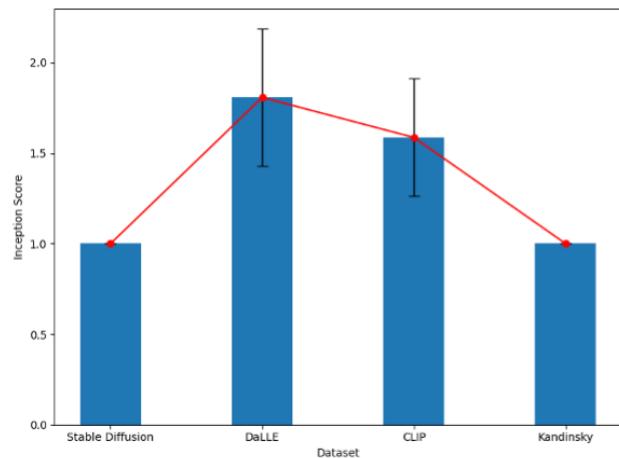


Figure 2. Generative Results across different models.

## 275 5.2. Evaluation

### 276 5.2.1. INCEPTION SCORE

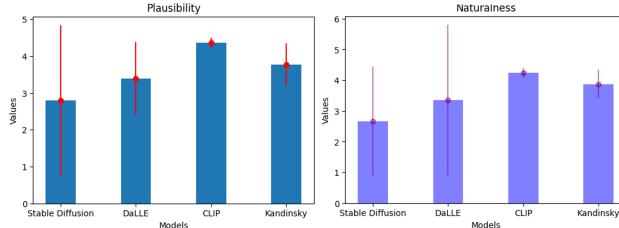


294 Figure 3. The empirical IS Score generated from the 4 models.

295  
296 Based on Figure 3, Stable Diffusion exhibits the lowest IS  
297 score. The IS score serves as a measure of a generative  
298 model's consistency in producing content. For instance, if  
299 the input text is "*the cat rides the horse*,"  $p_{\text{classifier}}(y|x_i(z))$   
300 represents the prediction of a classifier using a generated  
301 image  $x_i(z)$ , while  $\hat{p}(y, z)$  signifies the average distribution  
302 of all images generated from the model based on the  
303 input text "*the cat rides the horse*." Although these methods  
304 aim to offer a variety of images for users to choose from,  
305 consistency in content is still desired.

306 A small average KL distance indicates that all generated  
307 images from the model closely resemble the classifier's  
308 predictions. If the classifier is appropriately selected, it  
309 suggests that all images from the generative model belong  
310 to the same class as indicated by the classifier. Based on  
311 this understanding, we observe that both Stable Diffusion  
312 and Kandinsky's method notably outperform the remaining  
313 two methods.

### 314 5.2.2. HUMAN EVALUATION



325 Figure 4. Plausibility and Naturalness Evaluation.

326 The human evaluation revealed CLIP generated the most  
327 plausible images compared to the other models. Survey  
328

329 results showed most annotators felt CLIP and Kandinsky  
330 produced more sensible outputs than the rest. Although  
331 DaLLE achieved the highest IS score - indicating greater  
332 diversity - participants found its artistic style less realistic.  
333 Still, having more deviation bolstered diversity, playing to  
334 DaLLE's strength. In terms of naturalness though, DaLLE  
335 underperformed other models due to its creative stylization."

## 5.3. ControlNet

Control network is an extension of diffusion model, which shows a promising new method of adding spatial control to the process of T2I conversion by allowing additional conditional inputs (such as segmentation maps, depth maps, human bodies, gesture key points, etc.). Under the backbone of diffusion network, it provides a more accurate control direction in the stable diffusion backbone network than only obtaining clues from text. This allows AI to produce images that are more in line with our visual imagination, rather than merely fusions of features. A key innovation of ControlNet is its method of fine-tuning powerful reservation functions, which spread steadily in an end-to-end way without causing catastrophic forgetting or quality degradation. This makes it possible to learn specialized conditional tasks with smaller datasets. Collecting labeled elements and data is always expensive. Therefore, ControlNet limits the infinity of characteristics from the representative features of a small amount of data, giving us a great opportunity to get the desired results.

To assess the effectiveness of ControlNet, we conducted two tests on our sentence generation tasks:

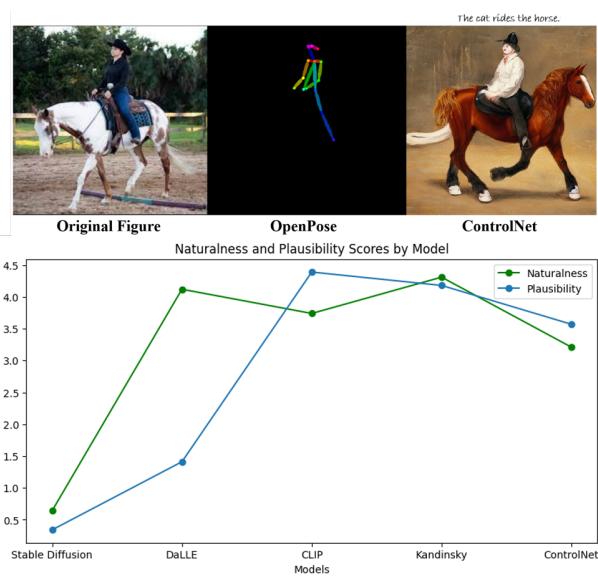
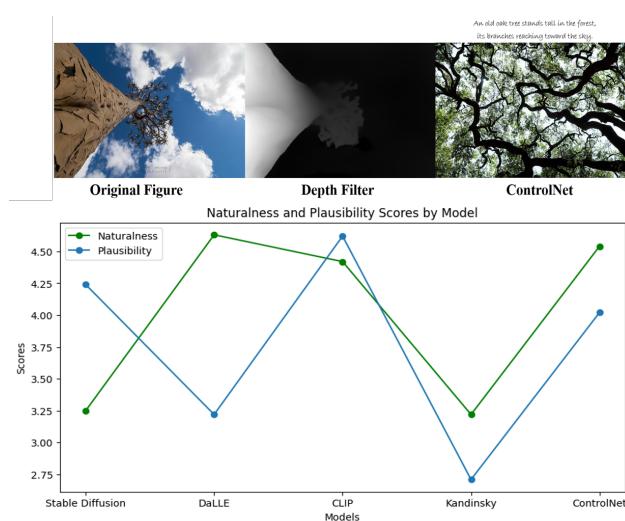


Figure 5. T2I Generation with Pose-strengthened ControlNet.

330 5.3.1. THE CAT RIDES THE HORSE  
331

332 For this task, we selected a relevant image of a cowgirl online  
333 and utilized the OpenPose toolkit to extract key body  
334 points representing her posture. These points were then  
335 input into ControlNet to constrain the movements of the  
336 main character described in the input text, in this case, the  
337 cat. As a result, 'riding' confined the cat's gesture, reflecting  
338 the expected posture. Evaluation of the images generated  
339 by ControlNet and other models revealed a significant  
340 improvement over the original Stable Diffusion method,  
341 approaching the quality of CLIP and Kandinsky. Human  
342 annotators consistently rated the naturalness and plausibility  
343 of ControlNet-generated images higher.

362 Figure 6. T2I Generation with Pose-strengthened ControlNet.  
363364 5.3.2. AN OLD OAK TREE STANDS TALL IN THE FOREST,  
365 ITS BRANCHES REACHING TOWARD THE SKY  
366

367 In this task, the phrase "reaching toward the sky" describes  
368 the dynamic vision of the disparity between the tree's  
369 branches and its trunk, highlighting the need to represent  
370 the tree's depth accurately. Based on the provided depth  
371 filter, ControlNet adeptly captures the depth and reach of  
372 tree branches through spatial and style constraints, resulting  
373 in highly realistic renderings. Comparative assessment  
374 against DaLLE and the original Diffusion model confirmed  
375 ControlNet's superior quality, as determined through qualitative  
376 human assessment. These results suggest that ControlNet  
377 successfully specifies the necessary spatial and style  
378 constraints through fine-tuning specific images. These experiments  
379 underscore ControlNet's robust performance in enhancing image generation tasks, contributing to generative  
380 models' advancements.

## Conclusion and Future Direction

In this study, we comprehensively evaluated four state-of-the-art text-to-image synthesis models: Stable Diffusion, DALL-E, CLIP, and Kandinsky. Our qualitative and quantitative comparisons revealed each model's strengths and limitations in generating high-quality, diverse, and semantically aligned images from textual descriptions. The results highlighted the superior performance of Controllable Generative AI, particularly when integrated with Large Language Models (LLMs) for enhanced biomimetic design processes. Future research should refine the proposed framework for Controllable Generative AI in Biomimetic Design, exploring novel techniques to incorporate real-world constraints and domain-specific knowledge into the creative process. Additionally, investigating the scalability and adaptability of these models to various application domains, such as architecture, product design, and biotechnology, will be crucial for unlocking their full potential. By bridging the gap between AI-generated creativity and practical feasibility, we can pave the way for innovative and sustainable solutions inspired by nature's wisdom.

## References

- Baltrušaitis, T., Ahuja, C., and Morency, L.-P. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018.
- Du, N., Huang, Y., Dai, A. M., Tong, S., Lepikhin, D., Xu, Y., Krikun, M., Zhou, Y., Yu, A. W., Firat, O., et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pp. 5547–5569. PMLR, 2022.
- Esser, P., Rombach, R., and Ommer, B. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12873–12883, 2021.
- Kang, M., Zhu, J.-Y., Zhang, R., Park, J., Shechtman, E., Paris, S., and Park, T. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10124–10134, 2023.
- Li, B. Z., Nye, M., and Andreas, J. Implicit representations of meaning in neural language models. *arXiv preprint arXiv:2106.00737*, 2021.
- Li, R., Li, W., Yang, Y., Wei, H., Jiang, J., and Bai, Q. Swinv2-imagen: Hierarchical vision transformer diffusion models for text-to-image generation. *Neural Computing and Applications*, pp. 1–16, 2023.

- 385 Liu, A. H., Jin, S., Lai, C.-I. J., Rouditchenko, A., Oliva,  
 386 A., and Glass, J. Cross-modal discrete representation  
 387 learning. *arXiv preprint arXiv:2106.05438*, 2021.
- 388 Mayerson, D. R. Gravitational multipoles in general stationary  
 389 spacetimes. *SciPost Physics*, 15(4):154, 2023.
- 390 Meng, C., Song, Y., Song, J., Wu, J., Zhu, J.-Y., and Ermon,  
 391 S. Sdedit: Image synthesis and editing with stochastic  
 392 differential equations. *arXiv preprint arXiv:2108.01073*,  
 393 2021.
- 394 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,  
 395 Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J.,  
 396 et al. Learning transferable visual models from natural  
 397 language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- 398 Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Rad-  
 399 ford, A., Chen, M., and Sutskever, I. Zero-shot text-to-  
 400 image generation. In *International conference on machine learning*, pp. 8821–8831. Pmlr, 2021.
- 401 Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M.  
 402 Hierarchical text-conditional image generation with clip  
 403 latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 404 Razzhigaev, A., Shakhmatov, A., Maltseva, A., Arkhip-  
 405 kin, V., Pavlov, I., Ryabov, I., Kuts, A., Panchenko, A.,  
 406 Kuznetsov, A., and Dimitrov, D. Kandinsky: an improved  
 407 text-to-image synthesis with image prior and latent diffu-  
 408 sion. *arXiv preprint arXiv:2310.03502*, 2023.
- 409 Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and  
 410 Ommer, B. High-resolution image synthesis with latent  
 411 diffusion models. In *Proceedings of the IEEE/CVF con-  
 412 ference on computer vision and pattern recognition*, pp.  
 413 10684–10695, 2022.
- 414 Sauer, A. and Geiger, A. Counterfactual generative net-  
 415 works. *arXiv preprint arXiv:2101.06046*, 2021.
- 416 Van Looveren, A., Klaise, J., Vacanti, G., and Cobb, O.  
 417 Conditional generative models for counterfactual expla-  
 418 nations. *arXiv preprint arXiv:2101.10123*, 2021.
- 419 Vincent, J. F., Bogatyreva, O. A., Bogatyrev, N. R., Bowyer,  
 420 A., and Pahl, A.-K. Biomimetics: its practice and theory.  
 421 *Journal of the Royal Society Interface*, 3(9):471–482,  
 422 2006.
- 423 Wang, Y., Qi, L., Chen, Y.-C., Zhang, X., and Jia, J. Image  
 424 synthesis via semantic composition. In *Proceedings of  
 425 the IEEE/CVF International Conference on Computer  
 426 Vision*, pp. 13749–13758, 2021.
- 427 Yadav, A. and Vishwakarma, D. K. Sentiment analysis  
 428 using deep learning architectures: a review. *Artificial  
 429 Intelligence Review*, 53(6):4335–4385, 2020.
- 430 Zameshina, M., Teytaud, O., and Najman, L. Diverse diffu-  
 431 sion: Enhancing image diversity in text-to-image genera-  
 432 tion. *arXiv preprint arXiv:2310.12583*, 2023.
- 433 Zhang, L., Rao, A., and Agrawala, M. Adding conditional  
 434 control to text-to-image diffusion models. In *Proceedings  
 435 of the IEEE/CVF International Conference on Computer  
 436 Vision*, pp. 3836–3847, 2023.
- 437 Zhou, Y. and Shimada, N. Vision+ language applications: A  
 438 survey. In *Proceedings of the IEEE/CVF Conference on  
 439 Computer Vision and Pattern Recognition*, pp. 826–842,  
 440 2023.

## 440 A. Appendix

### 441 A.1. Biomimetic Design: High-Speed Train

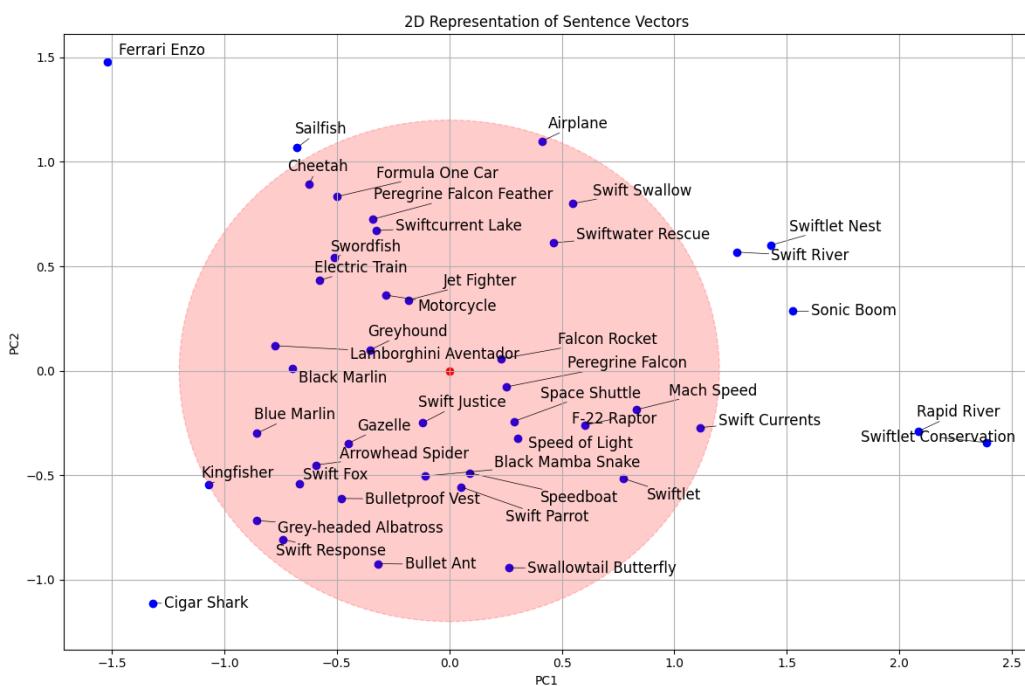
443 The generative models discussed previously have showcased remarkable creativity, enhancing the artistic domain of  
 444 Biomorphism and displaying versatility across various creative avenues. However, creativity alone may not suffice to  
 445 meet the demands of Biomimeticists. Drawing inspiration from nature's intricate designs and functionalities, biomedical  
 446 products often encounter challenges in transitioning from concept to reality. While Biomimetic scientists leverage their  
 447 understanding of nature to design and meet requirements, sourcing creatures with desired abilities can pose difficulties.  
 448 While our previous experiments highlighted ControlNet's potential to balance creativity and reliability, they only showcased  
 449 the model's generative capabilities. Here, we propose a framework that integrates the generalization ability of Large  
 450 Language Models (LLMs) with ControlNet to replicate a successful innovation, the High-Speed Train.

451 In this scenario, following the pipeline in Figure 1, the design process for the "High-Speed Train" unfolds as follows:

- 452 • **Entity Generation and Selection:** Generate a table of entities related to fast.

	Creature/Object	Description
0	Peregrine Falcon	Sleek pointed wings, Streamlined body long tai...
1	Cheetah	Low-slung body slender limbs, Muscular shoulde...
2	Sailfish	Torpedo-shaped body pointed bill, Muscular for...

- 453 • **Entity Projection and Outlier Removal:** Utilize Natural Language Processing (NLP) methods to identify and remove  
 454 outliers based on adjective similarity.



- 455 • **Adjective Collection and Filtering:** Compile all existing adjective terms and filter out those related to shape,  
 456 particularly addressing conflicting adjectives such as 'wide' and 'narrow,' 'tall' and 'short.'
- 457 • **Entity Selection using LLMs:** Leveraging OpenAI, we identify two entities: "Rocket" and "Aerodynamic Spacecraft."  
 458 We acquire an image of a rocket online and convert it into a canny image.

495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507508 ▼ Step 1. Pick up adjectives related to shape  
509  
510

```
'bulbous', 'angular', 'cylindrical', 'sleek', 'streamlined',
'pointed', 'lean', 'narrow', 'wide', 'deep', 'large',
'short', 'compact', 'broad', 'tall', 'small'
```

511 Step 2. Remove conflicting adjectives  
512  
513

```
'bulbous', 'angular', 'cylindrical', 'sleek', 'streamlined', 'pointed', 'lean', 'compact'
```

514 Step 3. Find objects that best fits the remaining adjectives:  
515  
516

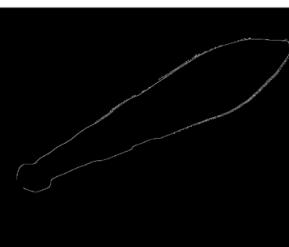
```
'Rocket', 'Aerodynamic Spacecraft'
```

- **Morphology-Constrained T2I Transformation:** With the constraint of the rocket's contour, the ControlNet generates a promising output for the text "A train runs fast."

517 A train runs fast



518 Rocket



519 Rocket's Contour



520 ControlNet

521 Based on the generated result of the ControlNet, it's evident that the contour of the rocket effectively shapes the train, giving  
 522 it a sleeker, more streamlined appearance with a pointed head. This design adjustment aids in reducing air resistance. Opting  
 523 for the contour of the rocket with canny edges, rather than emphasizing depth or adding intricate details to the rocket's  
 524 edges, may not sufficiently constrain the generation process. Our goal is to focus on the shape of the train rather than the  
 525 patterns on it. Providing excessive detail for the constraint could potentially mislead the ControlNet's generation process,  
 526 resulting in inaccurate feature fusion.

527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549